

Evaluating the Quality of RDF Data Sets on Common Vocabularies in the Social, Behavioral, and Economic Sciences

Thomas Hartmann¹, Benjamin Zapilko¹, Joachim Wackerow¹, and Kai Eckert²

¹ GESIS – Leibniz Institute for the Social Sciences, Germany

`{firstname.lastname}@gesis.org`,

² University of Mannheim, Germany

`kai@informatik.uni-mannheim.de`

Abstract. From 2012 to 2015 together with other Linked Data community members and experts from the social, behavioral, and economic sciences (*SBE*), we developed diverse vocabularies to represent SBE metadata and tabular data in RDF. The *DDI-RDF Discovery Vocabulary (DDI-RDF)* is designed to support the dissemination, management, and reuse of unit-record data, i.e., data about individuals, households, and businesses, collected in form of responses to studies and archived for research purposes. The *RDF Data Cube Vocabulary (QB)* is a W3C recommendation for expressing *data cubes*, i.e. multi-dimensional aggregate data and its metadata. *Physical Data Description (PHDD)* is a vocabulary to model data in rectangular format, i.e., tabular data. The data could either be represented in records with character-separated values (CSV) or fixed length. The *Simple Knowledge Organization System (SKOS)* is a vocabulary to build knowledge organization systems such as thesauri, classification schemes, and taxonomies. XKOS is a SKOS extension to describe formal statistical classifications.

To ensure high quality of and trust in both metadata and data, their representation in RDF must satisfy certain criteria - specified in terms of RDF constraints. In this paper, we evaluate the data quality of 15,694 data sets (4.26 billion triples) of research data for the social, behavioral, and economic sciences obtained from 33 SPARQL endpoints. We checked 115 constraints on three different and representative SBE vocabularies (DDI-RDF, QB, and SKOS) by means of the *RDF Validator*, a validation environment which is available at <http://purl.org/net/rdfval-demo>.

Keywords: RDF Validation, RDF Constraints, DDI-RDF Discovery Vocabulary, RDF Data Cube Vocabulary, Thesauri, SKOS, Linked Data, Semantic Web

1 Introduction

For constraint formulation and RDF data validation, several languages exist or are currently developed. *Shape Expressions (ShEx)*, *Resource Shapes (ReSh)*, *Description Set Profiles (DSP)*, *OWL 2*, the *SPARQL Inferencing Notation (SPIN)*,

and *SPARQL* are the six most promising and widely used constraint languages. OWL 2 is used as a constraint language under the closed-world and unique name assumptions. The W3C currently develops *SHACL*, an RDF vocabulary for describing RDF graph structures. With its direct support of validation via SPARQL, SPIN is very popular and certainly plays an important role for future developments in this field. It is particularly interesting as a means to validate arbitrary constraint languages by mapping them to SPARQL [4]. Yet, there is no clear favorite and none of the languages is able to meet all requirements raised by data practitioners. Further research and development therefore is needed.

In 2013, the W3C organized the RDF Validation Workshop,³ where experts from industry, government, and academia discussed first use cases for constraint formulation and RDF data validation. In 2014, two working groups on RDF validation have been established to develop a language to express constraints on RDF data: the *W3C RDF Data Shapes Working Group*⁴ (33 participants of 19 organizations) and the *DCMI RDF Application Profiles Task Group*⁵ (29 people of 22 organizations) which among others bundles the requirements of data institutions of the cultural heritage sector and the *social, behavioral, and economic (SBE) sciences* and represents them in the W3C group.

Within the DCMI task group, a collaboratively curated database of RDF validation requirements⁶ has been created which contains the findings of the working groups based on various case studies provided by data institutions [3]. It is publicly available and open for further contributions. The database connects requirements to use cases, case studies, and implementations and forms the basis of this paper. We distinguish 81 requirements to formulate constraints on RDF data; each of them corresponding to a constraint type.

We collected constraints for commonly used vocabularies in the SBE domain, either from the vocabularies themselves or from domain and data experts, in order to gain a better understanding about the role of certain requirements for data quality and to direct the further development of constraint languages. All in all, this lead to 115 constraints we implemented on three vocabularies. We let the experts classify the constraints according to the severity of their violation.

As we do not want to base our conclusions on the evaluation of vocabularies and constraint definitions alone, we conducted a large-scale experiment. For all these implemented 115 constraints, we evaluated the data quality of 15,694 data sets (4.26 billion triples) of SBE research data on three common vocabularies in SBE sciences (DDI-RDF, QB, SKOS) obtained from 33 SPARQL endpoints.

2 Common Vocabularies in SBE Sciences

We took all well-established and newly developed SBE vocabularies into account and defined constraints for three vocabularies commonly used in the SBE sciences

³ <http://www.w3.org/2012/12/rdf-val/>

⁴ <http://www.w3.org/2014/rds/charter>

⁵ <http://wiki.dublincore.org/index.php/RDF-Application-Profiles>

⁶ Online available at: <http://purl.org/net/rdf-validation>

which are briefly introduced in the following. We analyzed actual data according to constraint violations, as for these vocabularies large data sets are already published.

SBE sciences require high-quality data for their empirical research. For more than a decade, members of the SBE community have been developing and using a metadata standard, composed of almost twelve hundred metadata fields, known as the *Data Documentation Initiative (DDI)*,⁷ an XML format to disseminate, manage, and reuse data collected and archived for research [8]. In XML, the definition of schemas containing constraints and the validation of data according to these constraints is commonly used to ensure a certain level of data quality. With the rise of the Web of Data, data professionals and institutions are very interested in having their data be discovered and used by publishing their data directly in RDF or at least publish accurate metadata about their data to facilitate data integration. Therefore, not only established vocabularies like SKOS are used; recently, members of the SBE and Linked Data community developed with the *DDI-RDF Discovery Vocabulary (DDI-RDF)*⁸ a means to expose *DDI* metadata as Linked Data.

The data most often used in research within SBE sciences is *unit-record data*, i.e., data collected about individuals, businesses, and households, in form of responses to studies or taken from administrative registers such as hospital records, registers of births and deaths. A *study* represents the process by which a data set was generated or collected. The range of unit-record data is very broad - including census, education, health data and business, social, and labor force surveys. This type of research data is held within data archives or data libraries after it has been collected, so that it may be reused by future researchers. By its nature, unit-record data is highly confidential and access is often only permitted for qualified researchers who must apply for access. Researchers typically represent their results as aggregated data in form of multi-dimensional tables with only a few columns: so-called *variables* such as *sex* or *age*. Aggregated data, which answers particular research questions, is derived from unit-record data by statistics on groups or aggregates such as frequencies and arithmetic means. The purpose of publicly available aggregated data is to get a first overview and to gain an interest in further analyses on the underlying unit-record data. For more detailed analyses, researchers refer to unit-record data including additional variables needed to answer subsequent research questions.

Formal childcare is an example of an aggregated variable which captures the measured availability of childcare services in percent over the population in European Union member states by the dimensions *year*, *duration*, *age* of the child, and *country*. Variables are constructed out of values (of one or multiple datatypes) and/or code lists. The variable *age*, e.g., may be represented by values of the datatype *xsd:nonNegativeInteger* or by a code list of age clusters (e.g., '0 to 10' and '11 to 20'). The *RDF QB Vocabulary (QB)*⁹ is a W3C rec-

⁷ <http://www.ddialliance.org/Specification/>

⁸ <http://rdf-vocabulary.ddialliance.org/discovery.html>

⁹ <http://www.w3.org/TR/vocab-data-cube/>

ommendation for representing *QBs*, i.e., multi-dimensional aggregated data, in RDF [6]. A *qb:DataStructureDefinition* contains metadata of the data collection. The variable *formal childcare* is modeled as *qb:measure*, since it stands for what has been measured in the data collection. *Year*, *duration*, *age*, and *country* are *qb:dimensions*. Data values, i.e., the availability of childcare services in percent over the population, are collected in a *qb:DataSet*. Each data value is represented inside a *qb:Observation* which contains values for each dimension.

For more detailed analyses we refer to the underlying unit-record data. The aggregated variable *formal childcare* is calculated on the basis of six unit-record variables (i.a., *Education at pre-school*) for which detailed metadata is given (i.a., code lists) enabling researchers to replicate the results shown in aggregated data tables. *DDI-RDF* is used to represent metadata on unit-record data in RDF. The study (*disco:Study*) for which the unit-record data has been collected contains eight data sets (*disco:LogicalDataSet*) including variables (*disco:Variable*) like the six ones needed to calculate the variable *formal childcare*.

The *Simple Knowledge Organization System (SKOS)* is reused to a large extend to build SBE vocabularies. The codes of the variable *Education at pre-school* are modeled as *skos:Concepts* and a *skos:OrderedCollection* organizes them in a particular order within a *skos:memberList*. A variable may be associated with a theoretical concept (*skos:Concept*) and *skos:narrower* builds the hierarchy of theoretical concepts within a *skos:ConceptScheme* of a study. The variable *Education at pre-school* is assigned to the theoretical concept *Child Care* which is a narrower concept of the top concept *Education*. Controlled vocabularies (*skos:ConceptScheme*), serving as extension and reuse mechanism, organize types (*skos:Concept*) of descriptive statistics (*disco:SummaryStatistics*) like minimum, maximum, and arithmetic mean.

3 Classification of Constraint Types and Constraints

To gain better insights into the role that certain types of constraints play for the quality of RDF data, we use two simple classifications: on the one hand, we classify RDF constraint types whether they are expressible by different types of constraint languages and on the other hand, we classify constraints formulated for a given vocabulary according to the perceived severity of their violation.

Within the working groups, we identified by today 81 requirements to formulate RDF constraints (e.g., *R-75: minimum qualified cardinality restrictions*); each of them corresponding to an RDF constraint type.¹⁰ Within a technical report, we explain each requirement/constraint type in detail and give examples for each expressed by different constraint languages [5]. We provide mappings to representations in Description Logics (DL) [2] to logically underpin each requirement and to determine which DL constructs are needed to express each constraint type. For the three vocabularies, several SBE domain experts determined the default severity level of the 115 concrete constraints, which we pub-

¹⁰ Constraint types and constraints are uniquely identified by alphanumeric technical identifiers like *R-71-CONDITIONAL-PROPERTIES*

lished in a technical report [7]. In the following, we summarize the classifications of constraint types and constraints for the purpose of our evaluation.

3.1 Classification of Constraint Types according to the Expressivity of Constraint Languages

According to the expressivity of constraint languages, the complete set of constraint types encompasses three not disjoint sets of constraint types:

1. *RDFS/OWL Based*
2. *Constraint Language Based*
3. *SPARQL Based*

RDFS/OWL Based. The modeling languages RDFS and OWL are typically used to formally specify vocabularies. *RDFS/OWL Based* denotes the set of constraint types which can be formulated with RDFS/OWL axioms which we use in terms of constraints with CWA/UNA semantics and without reasoning.¹¹

Constraint Language Based. We further distinguish *Constraint Language Based* as the set of constraint types that can be expressed by common classical declarative high-level constraint languages like ShEx, ReSh, and DSP. There is a strong overlap between *RDFS/OWL* and *Constraint Language Based* constraint types as in many cases constraint types are expressible by both classical constraint languages and OWL. SPARQL, however, is considered as a low-level implementation language in this context. In contrast to SPARQL, high-level constraint languages are comparatively easy to understand and constraints can be formulated more concisely. Declarative languages may be placed on top of SPARQL when using it as an implementation language.

SPARQL Based. The set *SPARQL Based* encompasses constraint types that are not expressible by RDFS/OWL or common high-level constraint languages but by plain SPARQL.

3.2 Classification of RDF Constraints according to the Severity of Constraint Violations

A concrete constraint is instantiated from one of the 81 constraint types and is defined for a specific vocabulary. It does not make sense to determine the severity of constraint violations of an entire constraint type, as the severity depends on the individual context and vocabulary. SBE experts determined the default *severity level*¹² for each constraint to indicate how serious the violation of the constraint is. We use the classification system of log messages in software

¹¹ The entailment regime is to be decided by the implementers. It is our point that reasoning affects validation and that a proper definition of the reasoning to be applied is needed.

¹² The possibility to define severity levels in vocabularies is in itself a requirement (*R-158*).

development like *Apache Log4j 2* [1], the *Java Logging API*¹³ and the *Apache Commons Logging API*¹⁴ as many data practitioners also have experience in software development and software developers intuitively understand these levels. We simplify this commonly accepted classification system and distinguish the three severity levels (1) *informational*, (2) *warning*, and (3) *error*. Violations of *informational* constraints point to desirable but not necessary data improvements to achieve RDF representations which are ideal in terms of syntax and semantics of used vocabularies. *Warnings* are syntactic or semantic problems which typically should not lead to an abortion of data processing. *Errors*, in contrast, are syntactic or semantic errors which should cause the abortion of data processing. Although we provide default severity levels for each constraint, validation environments should enable users to adapt the severity levels of constraints according to their individual needs.

4 Evaluation

In this section, we describe our results based on an automatic constraint checking of a large data set. Despite the large volume of the data set in general, we have to keep in mind that this study only uses data for three vocabularies. As described in Section 2, for other vocabularies there is often not (yet) enough data openly available to draw general conclusions. The three vocabularies, however, are representative, cover different aspects of SBE data, and are also a mixture of widely adopted and accepted well-established vocabularies (QB, SKOS) and a vocabulary under development (DDI-RDF¹⁵).

4.1 Experimental Setup

On the three vocabularies (DDI-RDF, QB, SKOS), we identified and classified 115 constraints¹⁶ which we implemented for data validation. We ensured that the implementation of the constraints is equally distributed over the classes and vocabularies we have. We then evaluated the data quality of 15,694 data sets (4.26 billion triples) of SBE research data using these 115 constraints, obtained from 33 SPARQL endpoints.

Table 1 lists the number of validated data sets and the overall sizes in terms of triples for each of the vocabularies. We validated, i.a., (1) QB data sets published by the *Australian Bureau of Statistics*, the *European Central Bank*, and the *Organisation for Economic Co-operation and Development*, (2) SKOS thesauri like the *AGROVOC Multilingual agricultural thesaurus*, the *STW Thesaurus for Economics*, and the *Thesaurus for the Social Sciences*, and (3) DDI-RDF data

¹³ <http://docs.oracle.com/javase/7/docs/api/java/util/logging/Level.html>

¹⁴ <http://commons.apache.org/proper/commons-logging/>

¹⁵ Expected publication at the end of the year 2015

¹⁶ All 115 implemented constraints are online available at: <https://github.com/boschthomas/rdf-validation/tree/master/constraints>

sets provided by the *Microdata Information System*, the *Data Without Boundaries Discovery Portal*, the *Danish Data Archive*, and the *Swedish National Data Service*. We published the evaluation results for each QB data set in form of one document per SPARQL endpoint.¹⁷

Table 1: Validated Data Sets for each Vocabulary

Vocabulary Data Sets	Triples
QB	9,990 3,775,983,610
SKOS	4,178 477,737,281
DDI-RDF	1,526 9,673,055

Since the validation of each of the 81 constraint types can be implemented using SPARQL, we use *SPIN*, a SPARQL-based way to formulate and check constraints, as basis to develop a validation environment to validate RDF data according to constraints expressed by arbitrary constraint languages¹⁸ [4]. The *RDF Validator*¹⁹ can directly be used to validate arbitrary RDF data for the three vocabularies. Additionally, own constraints on any vocabulary can be defined using several constraint languages. The SPIN engine checks for each resource if it satisfies all constraints, which are associated with its assigned classes, and generates a result RDF graph containing information about all constraint violations. There is one SPIN construct template for each constraint type. A SPIN construct template contains a SPARQL CONSTRUCT query which generates constraint violation triples indicating the subject and the properties causing constraint violations and the reason why constraint violations have been raised. A SPIN construct template creates constraint violation triples if all triple patterns within the SPARQL WHERE clause match.

4.2 Evaluation Results

Tables 2 and 3 show the results of the evaluation, more specifically the constraints and the constraint violations, which are caused by these constraints, in percent; whereas the numbers in the first line indicate the absolute amount of constraints and violations. The constraints and their raised violations are grouped by vocabulary, which type of language the constraint types are formulated with, and their severity level. The numbers of validated triples and data sets differ between the vocabularies as we validated 3.8 billion QB, 480 million SKOS, and 10 million DDI-RDF triples. To be able to formulate findings which

¹⁷ Online available at: <https://github.com/boschthomas/rdf-validation/tree/master/evaluation/data-sets/data-cube>

¹⁸ Constraint language implementations online available at: <https://github.com/boschthomas/rdf-validation/tree/master/SPIN>

¹⁹ Online demo available at: <http://purl.org/net/rdfval-demo>, source code online available at: <https://github.com/boschthomas/rdf-validator>

apply for all vocabularies, we only use normalized relative values representing the percentage of constraints and violations belonging to the respective sets.

There is a strong overlap between *RDFS/OWL* and *Constraint Language Based* constraint types as in many cases constraint types are expressible by RDFS/OWL and classical constraint languages. This is the reason why the percentage values of constraints and violations grouped by the classification of constraint types according to the expressivity of constraint languages do not accumulate to 100%.

Table 2: Constraints and Constraint Violations (1)

	DDI-RDF		QB	
	C	CV	C	CV
	78	3,575,002	20	45,635,861
<i>SPARQL</i>	29.5	34.7	60.0	100.0
<i>CL</i>	64.1	65.3	40.0	0.0
<i>RDFS/OWL</i>	66.7	65.3	40.0	0.0
<i>info</i>	56.4	52.6	0.0	0.0
<i>warning</i>	11.5	29.4	15.0	99.8
<i>error</i>	32.1	18.0	85.0	0.3

C (constraints), *CV* (constraint violations)

Table 3: Constraints and Constraint Violations (2)

	SKOS		Total	
	C	CV	C	CV
	17	5,540,988	115	54,751,851
<i>SPARQL</i>	100.0	100.0	63.2	78.2
<i>CL</i>	0.0	0.0	34.7	21.8
<i>RDFS/OWL</i>	0.0	0.0	35.6	21.8
<i>info</i>	70.6	41.2	42.3	31.3
<i>warning</i>	29.4	58.8	18.7	62.7
<i>error</i>	0.0	0.0	39.0	6.1

C (constraints), *CV* (constraint violations)

4.3 Legend

In this sub-section, we describe how the tables in this paper should be read. Table 4 gives an overview over the symbols used in subsequent tables of the detailed evaluation.

Symbol	Description
✓	Validation Successful (without any constraint violation)
X	Constraint Violations
>X	Poor Performance/Scaling
X	Very Poor Performance/Scaling
(!)	Not Yet Implemented Constraint
(X)	The validation of X data sets could not be finished, due to SPARQL endpoints' technical restrictions (e.g., defined timeouts).
*	default severity level <i>informational</i>
**	default severity level <i>warning</i>
***	default severity level <i>error</i>

Table 4: Legend

- **Constraint Violations.** When constraints are violated, X indicates the number of raised constraint violation triples.
- **Poor Performance/Scaling.** The performance of the implementation of the underlying SPARQL CONSTRUCT query is too poor to get all resulting constraint violation triples. Therefore, a limit of X result constraint violation triples is set. It is likely that there are more than X constraint violations. Although the result set contains not the whole set of raised constraint violation triples, the constraint can be used as an indicator if there is data not conforming to the constraint and to resolve constraint violations step by step. As part of future work, the performance will be improved.
- **Very Poor Performance/Scaling.** The performance of the implementation of the underlying SPARQL CONSTRUCT query is too poor to get any results, even though a limit of result constraint violation triples is set. As part of future work, the performance will be improved.

5 Evaluation of Metadata on Unit-Record Data Sets (DDI-RDF)

In this section, the quality of the metadata on unit-record data sets (DDI-RDF) is evaluated by validating appropriate RDF constraints assigned to several RDF constraint types. First, we give an overview on the evaluated data sets and finally we provide details about the evaluation.

5.1 Data Sets Overview

Tables 5 and 7 give an overview on the evaluated DDI-RDF data sets, their abbreviations, and publicly available SPARQL endpoints. Table 6 comprehends

the number of triples, data sets, and instances of multiple vocabulary-specific classes.

Abbr.	DDI-RDF Data Sets
<i>Missy</i>	<i>Microdata Information System</i> ²⁰
<i>DwB</i>	<i>DwB Discovery Portal</i> ²¹
<i>DDA-SND</i>	<i>DDI-RDF</i> ²² provided by the <i>Danish Data Archive (DDA)</i> ²³ and <i>Swedish National Data Service (SND)</i> ²⁴

Table 5: DDI-RDF Data Sets Abbreviations

Data Sets	triples	Counts									
		disco:StudyGroup	disco:Study	disco:LogicalDataSet	disco:Universe	disco:Variable	disco:Question	disco:SummaryStatistics	disco:CategoryStatistics	skos:Concept	
<i>Missy</i>	5,068,838	6	45	159	1,125	21,040	0	0	0	147,193	
<i>DwB</i>	2,332,802	0	1,387	1,367	2,796	446,806	0	0	0	0	
<i>DDA-SND</i>	2,271,415	0	1,490	0	10,188	80,070	139,237	0	0	290,963	
Total	9,673,055			1,526							

Table 6: DDI-RDF Data Sets Overview

²⁰ <http://www.gesis.org/missy/eu/missy-home>

²¹ <http://dwb-dev.nsd.uib.no/portal>

²² <http://ddi-rdf.borsna.se/>

²³ <http://samfund.dda.dk/dda/default-en.asp>

²⁴ <http://snd.gu.se/en>

Data Sets	SPARQL Endpoint
<i>Missy</i>	http://svko-missy:8181/openrdf-workbench/repositories/native-java-store/summary
<i>DwB</i>	http://dwb-dev.nsd.uib.no/sparql
<i>DDA-SND</i>	http://ddi-rdf.borsna.se/endpoint/

Table 7: DDI-RDF SPARQL Endpoints

5.2 Detailed Evaluation

In this sub-section, we give details about the evaluation in form of diverse tables containing the number of constraint violations per evaluated data set and constraint of particular constraint types.

	Data Sets		
	<i>Missy</i>	<i>DwB</i>	<i>DDA-SND</i>
Existential Quantifications (1)			
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-01</i> ^{***}	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-02</i> ^{***}	7	17	1,490
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-03</i> [*]	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-04</i> [*]	11,021	445,381	62,260
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-05</i> [*]	✓	✓	139,237
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-06</i> [*]	12	1,367	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-07</i> [*]	6	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-08</i> [*]	45	1,387	1,490
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-09</i> [*]	6	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-10</i> [*]	45	1,387	1,490

Table 8: Evaluation of DDI-RDF Data Sets - Existential Quantifications (1)

Data Sets			
		Missy	DwB
DDA-SND			
Existential Quantifications (2)			
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-11</i> [*]	6	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-12</i> [*]	6	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-13</i> [*]	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-14</i> [*]	45	1,387	1,490
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-15</i> [*]	45	1,387	1,490
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-16</i> [*]	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-17</i> [*]	159	1,367	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-18</i> [*]	159	1,367	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-19</i> [*]	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-20</i> [*]	✓	1,367	✓

Table 9: Evaluation of DDI-RDF Data Sets - Existential Quantifications (2)

Data Sets			
	Missy	DwB	DDA-SND
Existential Quantifications (3)			
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-21</i> [*]	✓	1,367	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-22</i> [*]	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-23</i> [*]	6	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-24</i> [*]	45	1,387	1,490
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-25</i> [*]	45	1,387	1,490
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-26</i> [*]	45	1,387	1,490
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-27</i> ^{***}	✓	130	1,490
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-28</i> ^{**}	159	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-29</i> ^{**}	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-30</i> ^{**}	✓	✓	✓

Table 10: Evaluation of DDI-RDF Data Sets - Existential Quantifications (3)

Data Sets			
	Missy	DwB	DDA-SND
Existential Quantifications (4)			
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-31</i> **	159	1,367	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-32</i> ***	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-33</i> ***	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-34</i> ***	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-35</i> ***	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-36</i> ***	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-37</i> *	18,625	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-38</i> *	✓	✓	750
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-39</i> ***	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-40</i> *	✓	✓	139,237

Table 11: Evaluation of DDI-RDF Data Sets - Existential Quantifications (4)

Data Sets			
	Missy	DwB	DDA-SND
Existential Quantifications (5)			
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-41</i> *	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-42</i> *	✓	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-43</i> *	15,733	446,806	80,070
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-44</i> *	159	✓	✓
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-45</i> *	6,784	446,806	19,221
<i>DISCO-C-EXISTENTIAL-QUANTIFICATIONS-46</i> **	11,550	446,806	10,451

Table 12: Evaluation of DDI-RDF Data Sets - Existential Quantifications (5)

Data Sets			
Conditional Properties	<i>Missy</i>	<i>DwB</i>	<i>DDA-SND</i>
<i>DISCO-C-CONDITIONAL-PROPERTIES-01</i> ***	✓	✓	80,070
<i>DISCO-C-CONDITIONAL-PROPERTIES-02</i> **	12	✓	✓
<i>DISCO-C-CONDITIONAL-PROPERTIES-03</i> **	90	✓	2,980
<i>DISCO-C-CONDITIONAL-PROPERTIES-04</i> ***	6	✓	✓
<i>DISCO-C-CONDITIONAL-PROPERTIES-05</i> ***	45	1,387	1,490
<i>DISCO-C-CONDITIONAL-PROPERTIES-06</i> ***	✓	✓	✓

Table 13: Evaluation of DDI-RDF Data Sets - Conditional Properties

Data Sets			
Provenance	<i>Missy</i>	<i>DwB</i>	<i>DDA-SND</i>
<i>DISCO-C-PROVENANCE-01</i> *	6	✓	✓
<i>DISCO-C-PROVENANCE-02</i> *	45	1,387	1,490
<i>DISCO-C-PROVENANCE-03</i> *	159	1,367	✓
<i>DISCO-C-PROVENANCE-04</i> *	✓	1,367	✓

Table 14: Evaluation of DDI-RDF Data Sets - Provenance

Data Sets			
	<i>Missy</i>	<i>DwB</i>	<i>DDA-SND</i>
Labeling and Documentation			
<i>DISCO-C-LABELING-AND-DOCUMENTATION-01</i> [*]	6	✓	✓
<i>DISCO-C-LABELING-AND-DOCUMENTATION-02</i> [*]	45	1,387	1,490
<i>DISCO-C-LABELING-AND-DOCUMENTATION-03</i> [*]	159	1,367	✓
<i>DISCO-C-LABELING-AND-DOCUMENTATION-04</i> [*]	✓	1,367	✓
<i>DISCO-C-LABELING-AND-DOCUMENTATION-05</i> [*]	✓	✓	✓
<i>DISCO-C-LABELING-AND-DOCUMENTATION-06</i> [*]	21,040	446,806	80,070

Table 15: Evaluation of DDI-RDF Data Sets - Labeling and Documentation

Data Sets			
	<i>Missy</i>	<i>DwB</i>	<i>DDA-SND</i>
Data Model Consistency			
<i>DISCO-C-DATA-MODEL-CONSISTENCY-01</i> (!) ^{***}			
<i>DISCO-C-DATA-MODEL-CONSISTENCY-02</i> (!) ^{***}			
<i>DISCO-C-DATA-MODEL-CONSISTENCY-03</i> (!) ^{***}			
<i>DISCO-C-DATA-MODEL-CONSISTENCY-04</i> (!) ^{***}			
<i>DISCO-C-DATA-MODEL-CONSISTENCY-05</i> ^{***}	✓	✓	✓
<i>DISCO-C-DATA-MODEL-CONSISTENCY-06</i> (!) ^{***}			
<i>DISCO-C-DATA-MODEL-CONSISTENCY-07</i> (!) ^{***}			

Table 16: Evaluation of DDI-RDF Data Sets - Data Model Consistency

Comparison	Data Sets		
	Missy	DwB	DDA-SND
<i>DISCO-C-COMPARISON-VARIABLES-01 (!)**</i>			
<i>DISCO-C-COMPARISON-VARIABLES-02***</i>	21,040	446,806	80,070
<i>DISCO-C-COMPARISON-VARIABLES-03 (!)***</i>			
<i>DISCO-C-COMPARISON-VARIABLES-04 * 18,625</i>	✓	✓	
<i>DISCO-C-COMPARISON-VARIABLES-05 *** 159</i>	✓	✓	

Table 17: Evaluation of DDI-RDF Data Sets - Comparison

Mathematical Operations	Data Sets		
	Missy	DwB	DDA-SND
<i>DISCO-C-MATHEMATICAL-OPERATIONS-01 (!)***</i>			
<i>DISCO-C-MATHEMATICAL-OPERATIONS-02 (!)***</i>			
<i>DISCO-C-MATHEMATICAL-OPERATIONS-03 (!)***</i>			
<i>DISCO-C-MATHEMATICAL-OPERATIONS-04 (!)***</i>			
<i>DISCO-C-MATHEMATICAL-OPERATIONS-05 (!)***</i>			

Table 18: Evaluation of DDI-RDF Data Sets - Mathematical Operations

Data Sets		
Language Tags	Missy	DwB
DDA-SND		
<i>DISCO-C-LANGUAGE-TAG-MATCHING-01 (!)*</i>		
<i>DISCO-C-LANGUAGE-TAG-CARDINALITY-01 (!)*</i>		
<i>DISCO-C-LANGUAGE-TAG-CARDINALITY-02 (!)*</i>		
<i>DISCO-C-LANGUAGE-TAG-CARDINALITY-03 (!)*</i>		

Table 19: Evaluation of DDI-RDF Data Sets - Language Tags

Data Sets		
Aggregation	Missy	DwB
DDA-SND		
<i>DISCO-C-AGGREGATION-01 (!)*</i>		
<i>DISCO-C-AGGREGATION-02 (!)*</i>		
<i>DISCO-C-AGGREGATION-03 (!)*</i>		
<i>DISCO-C-AGGREGATION-04 (!)*</i>		
<i>DISCO-C-AGGREGATION-05 (!)*</i>		
<i>DISCO-C-AGGREGATION-06 (!)*</i>		
<i>DISCO-C-AGGREGATION-07 (!)*</i>		

Table 20: Evaluation of DDI-RDF Data Sets - Aggregation

DDI-RDF Constraints	Data Sets		
	Missy	DwB	DDA-SND
<i>DISCO-C-ALLOWED-VALUES-01</i> ***	✓	✓	✓
<i>DISCO-C-LITERAL-RANGES-01</i> ***	✓	✓	✓
<i>DISCO-C-INVVERSE-FUNCTIONAL-PROPERTIES-01</i> ***	✓	✓	✓
<i>DISCO-C-INVVERSE-FUNCTIONAL-PROPERTIES-02</i> ***	✓	✓	✓
<i>DISCO-C-CLASS-SPECIFIC-PROPERTY-RANGE-01</i> ***	✓	✓	✓
<i>DISCO-C-MEMBERSHIP-IN-CONTROLLED-VOCABULARIES-01</i> ***	✓	✓	✗
<i>DISCO-C-LITERAL-VALUE-COMPARISON-01</i> ***	✓	1,299	✓
<i>DISCO-C-CONTEXT-SPECIFIC-VALID-PROPERTIES-01</i> *	21,038	✓	✓
<i>DISCO-C-DATA-PROPERTY-FACETS-01</i> **	✓	✓	✓
<i>DISCO-C-DATA-PROPERTY-FACETS-02</i> **	✓	✓	✓

Table 21: Evaluation of DDI-RDF Data Sets - DDI-RDF Constraints (1)

Data Sets			
DDI-RDF Constraints	Missy	DwB	DDA-SND
<i>DISCO-C-VALUE-IS-VALID-FOR-DATATYPE-01</i> ***	30	6,932	✓
<i>DISCO-C-VALUE-IS-VALID-FOR-DATATYPE-02</i> ***		✓	✓
<i>DISCO-C-SUBSUMPTION-01</i> (!) *** ^B			
<i>DISCO-C-CLASS-EQUIVALENCE-01</i> (!)*			
<i>DISCO-C-SUB-PROPERTIES-01</i> (!) ***			
<i>DISCO-C-PROPERTY-DOMAIN-01</i> (!) ***			
<i>DISCO-C-PROPERTY-RANGES-01</i> (!) ***			
<i>DISCO-C-VERSE-OBJECT-PROPERTIES-01</i> (!) ***			
<i>DISCO-C-VERSE-OBJECT-PROPERTIES-02</i> (!) ***			
<i>DISCO-C-VERSE-OBJECT-PROPERTIES-03</i> (!) ***			
<i>DISCO-C-DISJOINT-PROPERTIES-01</i> (!) ***			

Table 22: Evaluation of DDI-RDF Data Sets - DDI-RDF Constraints (2)

	Data Sets
DDI-RDF Constraints	
<i>DISCO-C-ASYMMETRIC-OBJECT-PROPERTIES-01</i> (!)***	<i>Missy</i>
<i>DISCO-C-IRREFLEXIVE-OBJECT-PROPERTIES-01</i> (!)***	<i>DwB</i>
<i>DISCO-C-CLASS-SPECIFIC-IRREFLEXIVE-OBJECT-PROPERTIES-01</i> (!)***	
<i>DISCO-C-CLASS-SPECIFIC-IRREFLEXIVE-OBJECT-PROPERTIES-02</i> (!)***	
<i>DISCO-C-DISJOINT-CLASSES-01</i> (!)***	
<i>DISCO-C-EQUIVALENT-PROPERTIES-01</i> (!)*	
<i>DISCO-C-LITERAL-PATTERN-MATCHING-01</i> (!)*	
<i>DISCO-C-DISJUNCTION-01</i> (!)***	
<i>DISCO-C-UNIVERSAL-QUANTIFICATIONS-01</i> (!)***	
<i>DISCO-C-MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01</i> (!)***	

Table 23: Evaluation of DDI-RDF Data Sets - DDI-RDF Constraints (3)

	Data Sets
DDI-RDF Constraints	
<i>DISCO-C-MAXIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01</i>	<i>Missy</i>
<i>DISCO-C-EXACT-QUALIFIED-CARDINALITY-RESTRICTIONS-01</i>	<i>DwB</i>
<i>DISCO-C-CONTEXT-SPECIFIC-EXCLUSIVE-OR-OF-PROPERTY-GROUPS-01</i>	<i>DDA-SND</i>
<i>DISCO-C-IRI-PATTERN-MATCHING-01</i>	(!)*
<i>DISCO-C-ORDERING-01</i>	(!)*
<i>DISCO-C-ORDERING-02</i>	(!)*
<i>DISCO-C-ORDERING-03</i>	(!)*
<i>DISCO-C-STRING-OPERATIONS-01</i>	(!)*
<i>DISCO-C-CONTEXT-SPECIFIC-VALID-CLASSES-01</i>	(!)*
<i>DISCO-C-CONTEXT-SPECIFIC-VALID-PROPERTIES-01</i>	(!)*

Table 24: Evaluation of DDI-RDF Data Sets - DDI-RDF Constraints (4)

Data Sets	
DDI-RDF Constraints	
<i>DISCO-C-DEFAULT-VALUES-01 (!)*</i>	<i>Missy</i>
<i>DISCO-C-WHITESPACE-HANDLING-01 (!)*</i>	<i>DwB</i>
<i>DISCO-C-HTML-HANDLING-01 (!)*</i>	
<i>DISCO-C-HTML-HANDLING-02 (!)*</i>	
<i>DISCO-C-RECOMMENDED-PROPERTIES-01 (!)*</i>	
<i>DISCO-C-HANDLE-RDF-COLLECTIONS-01 (!)*</i>	
<i>DISCO-C-HANDLE-RDF-COLLECTIONS-02 (!)*</i>	
<i>DISCO-C-USE-SUB-SUPER-RELATIONS-IN-VALIDATION-01 (!)*</i>	
<i>DISCO-C-USE-SUB-SUPER-RELATIONS-IN-VALIDATION-02 (!)*</i>	
<i>DISCO-C-STRUCTURE-01 (!)***</i>	

Table 25: Evaluation of DDI-RDF Data Sets - DDI-RDF Constraints (5)

Data Sets	
DDI-RDF Constraints	
<i>DISCO-C-VOCABULARY-01 (!)***</i>	<i>Missy</i>
<i>DISCO-C-HTTP-URI-SCHEME-VIOLATION (!)***</i>	<i>DwB</i>

Table 26: Evaluation of DDI-RDF Data Sets - DDI-RDF Constraints (6)

6 Evaluation of Metadata and Data of Aggregated Data Sets (QB)

In this section, the quality of the metadata on aggregated data (QB) data sets and of the data sets themselves is evaluated by validating appropriate RDF con-

straints assigned to several RDF constraint types. First, we give an overview on the evaluated data sets and finally we provide details about the evaluation.

6.1 Data Sets Overview

There are websites giving an overview on available QB data sets²⁵. Tables 27 and 29 give an overview on the evaluated QB data sets, their abbreviations, and publicly available SPARQL endpoints. Table 28 comprehends the number of triples, data sets, and instances of multiple vocabulary-specific classes.

Abbr.	QB Data Sets
<i>ECB</i>	<i>European Central Bank</i> ²⁶
<i>UIS</i>	<i>UNESCO Institute for Statistics</i> ²⁷
<i>IMF</i>	<i>International Monetary Fund</i> ²⁸
<i>BFS</i>	<i>Bundesamt für Statistik - Swiss Federal Statistics</i> ²⁹
<i>FAO</i>	<i>Food and Agriculture Organization of the United Nations</i> ³⁰
<i>WB</i>	<i>World Bank</i> ³¹
<i>FRB</i>	<i>Federal Reserve Board</i> ³²
<i>TI</i>	<i>Transparency International</i> ³³
<i>OECD</i>	<i>Organisation for Economic Co-operation and Development</i> ³⁴
<i>BIS</i>	<i>Bank for International Settlements</i> ³⁵
<i>ABS</i>	<i>Australian Bureau of Statistics</i> ³⁶
<i>IEEE-VIS</i>	<i>IEEE VIS Source Data</i>
<i>ACORN-SAT</i>	<i>Australian Climate Observations Reference Network - Surface Air Temperature Dataset</i>
<i>HDP</i>	<i>HealthData.gov Platform (HDP) on the Semantic Web</i>
<i>Eurostat</i>	<i>The Eurostat Linked Data</i> (SPARQL endpoint unavailable)
<i>Asturias</i>	<i>Nomenclator Asturias</i> (SPARQL endpoint unavailable!)
<i>ISTAT</i>	<i>ISTAT Immigration (LinkedOpenData.it)</i> (SPARQL endpoint unavailable)
<i>ICANE</i>	<i>Statistical Office of Cantabria (Instituto Cántabro de Estadística, ICANE)</i> (SPARQL endpoint unavailable)
<i>EE-2009</i>	<i>European Election Results 2009</i> (SPARQL endpoint unavailable)
<i>EU-B</i>	<i>Standard Eurobarometer</i> (SPARQL endpoint unavailable)
<i>ECB-S</i>	<i>European Central Bank Statistics (PublicData.eu)</i> (SPARQL endpoint unavailable)
<i>CPV-2008</i>	<i>Common Procurement Vocabulary (CPV) 2008</i> (SPARQL endpoint unavailable)
<i>CPV-2003</i>	<i>Common Procurement Vocabulary (CPV) 2003</i> (SPARQL endpoint unavailable)

Table 27: QB Data Sets Abbreviations

²⁵ <http://270a.info/>;

<http://ontologycentral.com/>

<http://datahub.io/de/dataset?tags=format-qb>;

²⁶ <http://www.ecb.europa.eu/home/html/index.en.html>
²⁷ <http://www.uis.unesco.org/Pages/default.aspx>
²⁸ <http://www.imf.org/external/index.htm>
²⁹ <http://www.bfs.admin.ch/>
³⁰ <http://www.fao.org/home/en/>
³¹ <http://www.worldbank.org/>
³² <http://www.federalreserve.gov/>
³³ <http://www.transparency.org/>
³⁴ <http://www.oecd.org/>
³⁵ <http://www.bis.org/>
³⁶ <http://abs.gov.au/>

Counts						
Data Sets	triples	qb:DataSet	qb:DataStructureDefinition	qb:Observation	qb:Slice	
<i>ECB</i>	468,899,474	55	46	>11,000,000	428,698	
<i>UIS</i>	10,400,534	5	5	1,437,651	0	
<i>IMF</i>	35,688,446	4	8	3,603,719	0	
<i>BFS</i>	1,533,743	0	0	8	0	
<i>FAO</i>	53,000,000	10	10	>7,100,000	0	
<i>WB</i>	174,006,552	9,466	59	>17,000,000	0	
<i>FRB</i>	185,266,900	49	98	>9,500,000	0	
<i>TI</i>	52,233	6	6	3,928	0	
<i>OECD</i>	304,995,160	136	140	>12,000,000	0	
<i>BIS</i>	54,197,482	6	12	3,606,466	47,914	
<i>ABS</i>	2,357,400,000	253	257	>11,000,000	0	
<i>IEEE-VIS</i>	19,935,340	0	0	1,350	0	
<i>ACORN-SAT</i>	98,381,319	0	4	0	0	
<i>HDP</i>	12,226,427	0	0	0	0	
Total	3,775,983,610	9,990				

Table 28: QB Data Sets Overview

Data Sets	SPARQL Endpoints
<i>ECB</i>	http://ecb.270a.info/sparql
<i>UIS</i>	http://uis.270a.info/sparql
<i>IMF</i>	http://imf.270a.info/sparql
<i>BFS</i>	http:// bfs.270a.info/sparql
<i>FAO</i>	http://fao.270a.info/sparql
<i>WB</i>	http://worldbank.270a.info/sparql
<i>FRB</i>	http://frb.270a.info/sparql
<i>TI</i>	http://transparency.270a.info/sparql
<i>OECD</i>	http://oeecd.270a.info/sparql
<i>BIS</i>	http://bis.270a.info/sparql
<i>ABS</i>	http://abs.270a.info/sparql
<i>ACORN-SAT</i>	http://lab.environment.data.gov.au/sparql
<i>HDP</i>	http://healthdata.tw.rpi.edu/sparql

Table 29: QB SPARQL Endpoints

6.2 Detailed Evaluation

In this sub-section, we give details about the evaluation in form of diverse tables containing the number of constraint violations per evaluated data set and constraint of particular constraint types.

Data Model Consistency	Data Sets					
	ECB	UIS	IMF	BFS	FAO	WB
DATA-MODEL-CONSISTENCY-01**	✓ (2)	✓	✓	✓	✓	✓
DATA-MODEL-CONSISTENCY-02***	✓ (2)	✓	✓	✓	✓	✓
DATA-MODEL-CONSISTENCY-03***	✓ (2)	✓	✓	✓	✓	✓
DATA-MODEL-CONSISTENCY-04***	✓ (6)	✓	✓	✓	✓	14,372
DATA-MODEL-CONSISTENCY-05**	1,198,352 (50)	✗	✗	✓	✗	16,175,814 (42)
DATA-MODEL-CONSISTENCY-06***	✓ (2)	✓	✓	✓	✓	✓
DATA-MODEL-CONSISTENCY-07***	✓ (9)	✓	99,091	✓	✓	✓ (1)
DATA-MODEL-CONSISTENCY-08***	✓ (2)	✓	✓	✓	✓	✓
DATA-MODEL-CONSISTENCY-09***	✓ (2)	✓	✓	✓	✓	✓
DATA-MODEL-CONSISTENCY-10*** (!)	-	-	-	-	-	-
DATA-MODEL-CONSISTENCY-11**	6,511 (10)	✓	✓	✓	✓	✓

Table 30: Evaluation of QB Data Sets - Data Model Consistency (1)

Data Sets												
Data Model Consistency												
<i>DATA-MODEL-CONSISTENCY-01</i> **	✓			✓								
<i>DATA-MODEL-CONSISTENCY-02</i> ***	✓			✓								
<i>DATA-MODEL-CONSISTENCY-03</i> ***	✓			✓								
<i>DATA-MODEL-CONSISTENCY-04</i> ***	✓			✓				✓ (6)				
<i>DATA-MODEL-CONSISTENCY-05</i> **	✓	21,142,838 (116)		✗	6,997,098 (246)							
<i>DATA-MODEL-CONSISTENCY-06</i> ***	✓			✓				✓				
<i>DATA-MODEL-CONSISTENCY-07</i> ***	✓			✓				✓ (8)				
<i>DATA-MODEL-CONSISTENCY-08</i> ***	✓			✓				✓				
<i>DATA-MODEL-CONSISTENCY-09</i> ***	✓			✓				✓				
<i>DATA-MODEL-CONSISTENCY-10</i> *** (!)	-		-	-				-		-	-	-
<i>DATA-MODEL-CONSISTENCY-11</i> **	✓		✓		✓			✓		✓	✓	✓

Table 31: Evaluation of QB Data Sets - Data Model Consistency (2)

Data Sets												
Existential Quantifications												
<i>EXISTENTIAL-QUANTIFICATIONS-01</i> ***	9	✓	11	7	8	77	8	9	7	8	7	✓
<i>EXISTENTIAL-QUANTIFICATIONS-02</i> ***	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
<i>EXISTENTIAL-QUANTIFICATIONS-03</i> ***	✓	✓	✓	✓	✓	✓	59	✓	6	✓	✓	✓
<i>EXISTENTIAL-QUANTIFICATIONS-04</i> ***	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 32: Evaluation of QB Data Sets - Existential Quantifications

Data Sets										
Cardinality Restrictions	<i>ECB</i>	<i>UIS</i>	<i>IMF</i>	<i>BFS</i>	<i>FAO</i>	<i>WB</i>	<i>FRB</i>	<i>TI</i>	<i>OECD</i>	<i>BIS</i>
<i>MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01 (!)***</i>	-	-	-	-	-	-	-	-	-	-
<i>MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-02***</i>	X	118	8	8	30	✓	30	✓	X	12
<i>MAXIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01***</i>	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
<i>EXACT-UNQUALIFIED-CARDINALITY-RESTRICTIONS-01***</i>	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
<i>EXACT-QUALIFIED-CARDINALITY-RESTRICTIONS-02***</i>	✓	✓	✓	✓	✓	1	✓	✓	✓	✓

Table 33: Evaluation of QB Data Sets - Cardinality Restrictions (1)

Data Sets				
Cardinality Restrictions	<i>ABS</i>	<i>IEEE-VIS</i>	<i>ACORN-SAT</i>	<i>HDP</i>
<i>MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01 (!)***</i>	-	-	-	-
<i>MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-02***</i>	X	1,350	✓	✓
<i>MAXIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01***</i>	✓	(2)	✓	✓
<i>EXACT-UNQUALIFIED-CARDINALITY-RESTRICTIONS-01***</i>	✓	✓	✓	✓
<i>EXACT-QUALIFIED-CARDINALITY-RESTRICTIONS-02***</i>	✓	✓	✓	✓

Table 34: Evaluation of QB Data Sets - Cardinality Restrictions (2)

Data Sets												
Structure	ECB	UIS	IMF	BFS	FAO	WB	FRB	TI	OECD	BIS	ABS	IEEE-VIS
STRUCTURE-01***	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
STRUCTURE-02***	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 35: Evaluation of QB Data Sets - Structure

Data Sets												
Constraints	ECB	UIS	IMF	BFS	FAO	WB	FRB	TI	OECD	BIS	ABS	IEEE-VIS
PROPERTY-DOMAIN-01 (!)***												
PROPERTY-RANGES-01 (!)***												
DISJOINT-PROPERTIES-01 (!)***												
DISJOINT-CLASSES-01 (!)***												
EQUIVALENT-PROPERTIES-01 (!)*												
UNIVERSAL-QUANTIFICATIONS-01 (!)***												
MEMBERSHIP-IN-CONTROLLED-VOCABULARIES-01 (!)***												
CONTEXT-SPECIFIC-VALID-CLASSES-01 (!)*												
CONTEXT-SPECIFIC-VALID-PROPERTIES-01 (!)*												
RECOMMENDED-PROPERTIES-01 (!)*												
VALUE-IS-VALID-FOR-DATATYPE-01 (!)***												
VOCABULARY-01 (!)***												

Table 36: Evaluation of QB Data Sets - Constraints (1)

Data Sets	
Constraints	
<i>HTTP-URI-SCHEME-VIOLATION</i> (!)***	
	<i>ECB</i>
	<i>UIS</i>
	<i>IMF</i>
	<i>BFS</i>
	<i>FAO</i>
	<i>WB</i>
	<i>FRB</i>
	<i>TI</i>
	<i>OECD</i>
	<i>BIS</i>
	<i>ABS</i>
	<i>IEEE-VIS</i>
	<i>ACORN-SAT</i>
	<i>HDP</i>

Table 37: Evaluation of QB Data Sets - Constraints (2)

7 Evaluation of Metadata on Thesauri (SKOS)

In this section, the quality of the metadata on thesauri (SKOS) is evaluated by validating appropriate RDF constraints assigned to several RDF constraint types. First, we give an overview on the evaluated thesauri and finally we provide details about the evaluation.

7.1 Data Sets Overview

There is a website giving an overview on available SKOS data sets³⁷ and another one giving an overview on available thesauri³⁸. Tables 38 and 40 give an overview on the evaluated thesauri, their abbreviations, and publicly available SPARQL endpoints. Table 39 comprehends the number of triples, data sets, and instances of multiple vocabulary-specific classes.

³⁷ <http://datahub.io/de/dataset?tags=format-skos>

³⁸ <http://datahub.io/de/dataset?tags=thesaurus>

Abbr.	Thesauri
<i>TheSoz</i>	<i>Thesaurus for the Social Sciences</i> ³⁹
<i>STW</i>	<i>Thesaurus for Economics</i> ⁴⁰
<i>AGROVOC</i>	<i>AGROVOC Multilingual agricultural thesaurus</i> ⁴¹
<i>UNESCO</i>	<i>UNESCO Thesaurus</i> ⁴²
<i>TGN</i>	<i>The Getty Thesaurus of Geographic Names</i> ⁴³
<i>EARTH</i>	<i>Environmental Applications Reference Thesaurus</i> ⁴⁴
<i>ODT</i>	<i>Open Data Thesaurus</i> ⁴⁵
<i>SLD</i>	<i>Spanish Linguistic Datasets</i> ⁴⁶
<i>SSWT</i>	<i>Social Semantic Web Thesaurus</i> ⁴⁷
<i>GBA-GU</i>	<i>Thesaurus of the Geological Survey of Austria (GBA) - Geology Unit</i> ⁴⁸
<i>GBA-GTS</i>	<i>Thesaurus of the Geological Survey of Austria (GBA) - Geologic Time Scale</i> ⁴⁹
<i>GBA-L</i>	<i>Thesaurus of the Geological Survey of Austria (GBA) - Lithology</i> ⁵⁰
<i>GBA-LU</i>	<i>Thesaurus of the Geological Survey of Austria (GBA) - Lithotectonic Unit</i> ⁵¹
<i>GEMET</i>	<i>GEneral Multilingual Environmental Thesaurus</i> ⁵²
<i>EuroVoc</i>	<i>EuroVoc</i> ⁵³
<i>CECCT</i>	<i>Clean Energy and Climate Change Thesaurus</i> ⁵⁴

Table 38: Thesauri Abbreviations

³⁹ <http://www.ecb.europa.eu/home/html/index.en.html>

⁴⁰ <http://zbw.eu/stw/versions/latest/about>

⁴¹ <http://202.45.139.84:10035/catalogs/fao/repositories/agrovoc>

⁴² <http://skos.um.es/sparql/>

⁴³ <http://vocab.getty.edu/sparql>

⁴⁴ <http://linkeddata.ge.imati.cnr.it/resource/EARTH/>

⁴⁵ <http://vocabulary.semantic-web.at/PoolParty/wiki/OpenData>

⁴⁶ <http://linguistic.linkeddata.es>

⁴⁷ <http://vocabulary.semantic-web.at/PoolParty/wiki/semweb>

⁴⁸ <http://resource.geolba.ac.at/>

⁴⁹ <http://resource.geolba.ac.at/>

⁵⁰ <http://resource.geolba.ac.at/>

⁵¹ <http://resource.geolba.ac.at/>

⁵² <http://www.eionet.europa.eu/gemet/>

⁵³ <http://open-data.europa.eu/de/data/dataset/eurovoc>

⁵⁴ <http://data.reegle.info/thesaurus/guide>

Thesauri	triples	Counts					
		skos:Concept	skos:broader	skos:narrower	skos:hasTopConcept	skos:inScheme	
<i>TheSoz</i>	439,153	1	8,426	13,705	13,706	0	48,529
<i>STW</i>	221,668	1	13,468	13,732	13732	7	13,180
<i>AGROVOC</i>	6,080,477	1	32,310	33,507	33,507	25	32,310
<i>UNESCO</i>	288,346	9	26,714	20,028	20,028	607	32,009
<i>TGN</i>	16,112,321	8	2,898,775	0	0	0	1,453,767
<i>EARTH</i>	9,287,364	11	295,375	288,208	93,827	479	295,376
<i>ODT</i>	3,290	6	108	93	93	30	0
<i>SLD</i>	7,629,211	0	31,195	0	0	0	0
<i>SSWT</i>	64,698	9	2,127	2,300	2,301	38	0
<i>GBA-GU</i>	25,718	3	878	1,005	1,005	14	0
<i>GBA-GTS</i>	7,875	3	213	208	208	5	0
<i>GBA-L</i>	9,317	1	249	249	249	4	0
<i>GBA-LU</i>	9,504	3	364	359	359	7	0
<i>GEMET</i>	372,889,229	3,680	414,659	62,193	21,685	30,806	409,290
<i>EuroVoc</i>	64,477,774	439	79,557	6,922	0	532	14,428
<i>CECCT</i>	191,336	3	3,419	3,761	3,762	28	0
Total	477,737,281	4,178					

Table 39: Thesauri Overview

Thesauri	SPARQL Endpoints
<i>TheSoz</i>	http://lod.gesis.org/thesoz/sparql
<i>STW</i>	http://zbw.eu/beta/sparql/stw/query
<i>AGROVOC</i>	http://202.45.139.84:10035/catalogs/fao/repositories/agrovoc
<i>UNESCO</i>	http://skos.um.es/sparql/
<i>TGN</i>	http://vocab.getty.edu/
<i>EARTH</i>	http://linkeddata.ge.imati.cnr.it:8890/sparql
<i>ODT</i>	http://vocabulary.semantic-web.at/PoolParty/sparql/OpenData
<i>SLD</i>	http://linguistic.linkeddata.es/sparql
<i>SSWT</i>	http://vocabulary.semantic-web.at/PoolParty/sparql/semweb
<i>GBA-GU</i>	http://resource.geolba.ac.at/PoolParty/sparql/GeologicUnit
<i>GBA-GTS</i>	http://resource.geolba.ac.at/PoolParty/sparql/GeologicTimeScale
<i>GBA-L</i>	http://resource.geolba.ac.at/PoolParty/sparql/lithology
<i>GBA-LU</i>	http://resource.geolba.ac.at/PoolParty/sparql/tectonicunit
<i>GEMET</i>	http://semantic.eea.europa.eu/sparql
<i>Euro Voc</i>	http://open-data.europa.eu/de/linked-data
<i>CECCT</i>	http://poolparty.reegle.info/PoolParty/sparql/glossary

Table 40: Thesauri SPARQL Endpoints

7.2 Detailed Evaluation

In this sub-section, we give details about the evaluation in form of diverse tables containing the number of constraint violations per evaluated data set and constraint of particular constraint types.

Data Model Consistency		Data Sets									
<i>TheSoz</i>		<i>STW</i>	<i>AGROVOC</i>	<i>TGN</i>	<i>UNESCO</i>	<i>ODT</i>	<i>SSWT</i>	<i>GBA-GU</i>	<i>GBA-GTS</i>	<i>GBA-L</i>	<i>GBA-LU</i>
<i>DATA-MODEL-CONSISTENCY-01</i> (!)*											
<i>DATA-MODEL-CONSISTENCY-02</i> (!)*											
<i>DATA-MODEL-CONSISTENCY-03</i> (!)*											

Table 41: Thesauri Evaluation - Data Model Consistency (1)

Data Sets									
Data Model Consistency									
<i>DATA-MODEL-CONSISTENCY-01</i> (!)*									
<i>DATA-MODEL-CONSISTENCY-02</i> (!)*									
<i>DATA-MODEL-CONSISTENCY-03</i> (!)*									

Table 42: Thesauri Evaluation - Data Model Consistency (2)

Data Sets									
Labeling and Documentation									
<i>LABELING-AND-DOCUMENTATION-01</i> * 8,426	11,508	19,829	1,110	✗	36	1,475	5	2	✓ 107 486
<i>LABELING-AND-DOCUMENTATION-02</i> * >1	✗	>100	287	✗ ✓	✓	✓ ✓ ✓ ✓ ✓			
<i>LABELING-AND-DOCUMENTATION-03</i> * ✓	✓	1	14,114	✗ ✓	✓	1 ✓ ✓ ✓	1		
<i>LABELING-AND-DOCUMENTATION-04</i> (!)* ✓	✓	4	✓	1 2	2	1 ✓ ✓ ✓	✓	✓	7
<i>LABELING-AND-DOCUMENTATION-05</i> * 975,340	✓	✓	2	✓ ✓	✓	✓ ✓ ✓ ✓ ✓	✓	✓	
<i>LABELING-AND-DOCUMENTATION-06</i> * 975,340	✓	✓	2	✓ ✓	✓	✓ ✓ ✓ ✓ ✓	✓	✓	

Table 43: Thesauri Evaluation - Labeling and Documentation (1)

Data Sets							
Labeling and Documentation							
<i>LABELING-AND-DOCUMENTATION-01</i> *		264,687	X	54,911	31,195		
<i>LABELING-AND-DOCUMENTATION-02</i> *			X	X	X	✓	
<i>LABELING-AND-DOCUMENTATION-03</i> *		2	X	55,556	31,195		
<i>LABELING-AND-DOCUMENTATION-04 (!)</i> *							
<i>LABELING-AND-DOCUMENTATION-05</i> *		39	X	X	978		
<i>LABELING-AND-DOCUMENTATION-06</i> *	302	46,718	✓	✓			

Table 44: Thesauri Evaluation - Labeling and Documentation (2)

Data Sets										
Structure	<i>TheSoz</i>	<i>STW</i>	<i>AGROVOC</i>	<i>TGN</i>	<i>UNESCO</i>	<i>ODT</i>	<i>SSWT</i>	<i>GBA-GU</i>	<i>GBA-GTS</i>	<i>CECCT</i>
<i>STRUCTURE-01</i> **	1	1,074	✓	✓	1	5	1	✓	✓	✓
<i>STRUCTURE-02 (!)</i> *										
<i>STRUCTURE-03</i> **	✓	✓	✓	✓	84	✓	✓	✓	✓	✓
<i>STRUCTURE-04</i> *	2,906	8,046	726	✓	3,840	12	124	84	256	68
<i>STRUCTURE-05</i> *	✓	✓	✓	✓	X	90	5,150	✓	✓	✓
<i>STRUCTURE-06</i> *	1,457	37	✓	✓	X	✓	4	1	1	64
<i>STRUCTURE-07</i> **	40	5,370	✓	✓	X	✓	✓	✓	✓	✓
<i>STRUCTURE-08 (!)</i> ***										
<i>STRUCTURE-09</i> *	7,897	19,844	99	✓	552	2	16	26	✓	✓
<i>STRUCTURE-10</i> **	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 45: Thesauri Evaluation - Structure (1)

Data Sets						
Structure	EARTh		GEMET		EuroVoc	
						SLD
STRUCTURE-01 ^{**}	18,240	X	55,757	31,195		
STRUCTURE-02 (!) [*]						
STRUCTURE-03 ^{**}	39	4,244	✓	✓		
STRUCTURE-04 [*]	11,286	74	✓	✓		
STRUCTURE-05 [*]	✓	X	✓	✓		
STRUCTURE-06 [*]	239,346	X	13,876	✓		
STRUCTURE-07 ^{**}	110,015	X	366,155	155,975		
STRUCTURE-08 (!) ^{***}						
STRUCTURE-09 [*]	107,195	32	✓	✓		
STRUCTURE-10 ^{**}	27	2,122	✓	✓		

Table 46: Thesauri Evaluation - Structure (2)

Data Sets										
Language Tag Cardinality	<i>TheSoz</i>	<i>STW</i>	<i>AGROVOC</i>	<i>TGN</i>	<i>UNESCO</i>	<i>ODT</i>	<i>SSWT</i>	<i>GBA-GU</i>	<i>GBA-GTS</i>	<i>GBA-L</i>
LANGUAGE-TAG-CARDINALITY-01 ^{**}	9,435	13,468	98,894	✓	541	10,147	5,117	2,061	1,742	2,272
LANGUAGE-TAG-CARDINALITY-02 [*]	8,222	36,936	X	✓	265	3,627	2,212	635	631	1,253
LANGUAGE-TAG-CARDINALITY-03 [*]	8,222	✓	135	✓	✓	✓	✓	✓	✓	✓
LANGUAGE-TAG-CARDINALITY-04 [*]	✓	476	X	50	✓	✓	✓	✓	✓	✓

Table 47: Thesauri Evaluation - Language Tag Cardinality (1)

		Data Sets			
Language Tag	Cardinality	<i>EARTH</i>	<i>GEMET</i>	<i>EuroVoc</i>	<i>SLD</i>
<i>LANGUAGE-TAG-CARDINALITY-01</i> **	✗	2,318,895	✗	30,781	
<i>LANGUAGE-TAG-CARDINALITY-02</i> *	✗		✗	✗	✗
<i>LANGUAGE-TAG-CARDINALITY-03</i> *	224,206		✗	✗	31,195
<i>LANGUAGE-TAG-CARDINALITY-04</i> *	✗	✗	✓	✓	

Table 48: Thesauri Evaluation - Language Tag Cardinality (2)

		Data Sets											
Constraints		<i>TheSoz</i>	<i>STW</i>	<i>AGROVOC</i>	<i>TGN</i>	<i>UNESCO</i>	<i>ODT</i>	<i>SSWT</i>	<i>GBA-GU</i>	<i>GBA-GTS</i>	<i>GBA-L</i>	<i>GBA-LU</i>	<i>CECCT</i>
<i>PROPERTY-DOMAIN-01</i> (!)***													
<i>PROPERTY-RANGES-01</i> (!)***													
<i>DISJOINT-PROPERTIES-01</i> (!)***													
<i>DISJOINT-PROPERTIES-02</i> (!)***													
<i>DISJOINT-CLASSES-01</i> (!)***													
<i>EQUIVALENT-PROPERTIES-01</i> (!)*													
<i>UNIVERSAL-QUANTIFICATIONS-01</i> (!)***													
<i>CONTEXT-SPECIFIC-VALID-CLASSES-01</i> (!)*													
<i>CONTEXT-SPECIFIC-VALID-PROPERTIES-01</i> (!)*													
<i>RECOMMENDED-PROPERTIES-01</i> (!)*													
<i>VOCABULARY-01</i> (!)***													
<i>HTTP-URI-SCHEME-VIOLATION</i> (!)***													

Table 49: Thesauri Evaluation - Constraints (1)

Data Sets	
Constraints	
<i>PROPERTY-DOMAIN-01</i> (!)***	<i>EARTH</i>
<i>PROPERTY-RANGES-01</i> (!)***	<i>GEMET</i>
<i>DISJOINT-PROPERTIES-01</i> (!)***	<i>EuroVoc</i>
<i>DISJOINT-PROPERTIES-02</i> (!)***	
<i>DISJOINT-CLASSES-01</i> (!)***	<i>SLD</i>
<i>EQUIVALENT-PROPERTIES-01</i> (!)*	
<i>UNIVERSAL-QUANTIFICATIONS-01</i> (!)***	
<i>CONTEXT-SPECIFIC-VALID-CLASSES-01</i> (!)*	
<i>CONTEXT-SPECIFIC-VALID-PROPERTIES-01</i> (!)*	
<i>RECOMMENDED-PROPERTIES-01</i> (!)*	
<i>VOCABULARY-01</i> (!)***	
<i>HTTP-URI-SCHEME-VIOLATION</i> (!)***	

Table 50: Thesauri Evaluation - Constraints (2)

8 Evaluation of Metadata on Statistical Classifications (XKOS)

As part of future work, the quality of metadata on statistical classifications (XKOS) data sets will be evaluated by validating appropriate RDF constraints assigned to several RDF constraint types.

8.1 Data Sets Overview

Abbr.	Statistical Classifications
<i>NAF</i>	<i>Nomenclature d'activités française</i> ⁵⁵
<i>PCS</i>	<i>Nomenclature des Professions et Catégories Socioprofessionnelles</i> ⁵⁶
<i>CJ</i>	<i>Nomenclature des catégories juridiques</i> ⁵⁷
<i>ISIC</i>	
<i>ISCO</i>	<i>International Standard Classification of Occupations</i>

Table 51: Statistical Classifications Abbreviations

Nomenclature d'activités française (NAF) is the French refinement of the *NACE* classification expressed in XKOS having explanatory notes. *Nomenclature des Professions et Catégories Socioprofessionnelles (PCS)* and *Nomenclature des catégories juridiques (CJ)* are French classifications expressed in XKOS. The statistical classification *ISIC* has explanatory notes too.

9 Conclusion

We identified and published by today 81 types of constraints that are required by various stakeholders for data applications. In close collaboration with several domain experts for the social, behavioral, and economic sciences (SBE), we formulated and implemented 115 constraints on three different vocabularies (DDI-RDF, QB, and SKOS) and classified them according to their severity level and whether their type is expressible by different types of constraint languages - RDFS/OWL, high-level constraint languages, and SPARQL. Using these constraints, we evaluated the data quality of 15,694 data sets (4.26 billion triples) of research data for the SBE sciences obtained from 33 SPARQL endpoints.

References

1. Apache Software Foundation. Apache Log4j 2 v. 2.3 User's Guide. Technical report, Apache Software Foundation, May 2015. <http://logging.apache.org/log4j/2.x/log4j-users-guide.pdf>.
2. Franz Baader, Diego Calvanese, Deborah L. McGuinness, Daniele Nardi, and Peter F. Patel-Schneider, editors. *The Description Logic Handbook: Theory, Implementation, and Applications*. Cambridge University Press, New York, NY, USA, 2003.
3. Thomas Bosch and Kai Eckert. Requirements on RDF Constraint Formulation and Validation. In *Proceedings of the 14th DCMI International Conference on Dublin Core and Metadata Applications (DC 2014)*, Austin, Texas, USA, 2014. <http://dcevents.dublincore.org/IntConf/dc-2014/paper/view/257>.
4. Thomas Bosch and Kai Eckert. Towards Description Set Profiles for RDF using SPARQL as Intermediate Language. In *Proceedings of the 14th DCMI International Conference on Dublin Core and Metadata Applications (DC 2014)*, Austin, Texas, USA, 2014. <http://dcevents.dublincore.org/IntConf/dc-2014/paper/view/270>.
5. Thomas Bosch, Andreas Nolle, Erman Acar, and Kai Eckert. RDF Validation Requirements - Evaluation and Logical Underpinning. *Computing Research Repository (CoRR)*, abs/1501.03933, 2015. <http://arxiv.org/abs/1501.03933>.
6. Richard Cyganiak, Simon Field, Arofan Gregory, Wolfgang Halb, and Jeni Tennison. Semantic Statistics: Bringing Together SDMX and SCODO. In Christian Bizer, Tom Heath, Tim Berners-Lee, and Michael Hausenblas, editors, *Proceedings of the International World Wide Web Conference (WWW 2010), Workshop on Linked Data on the Web*, volume 628 of *CEUR Workshop Proceedings*, 2010. <http://ceur-ws.org/Vol-628/lidow2010-paper03.pdf>.

⁵⁵ <http://rdf.insee.fr/codes/index.html>

⁵⁶ <http://rdf.insee.fr/codes/index.html>

⁵⁷ <http://rdf.insee.fr/codes/index.html>

7. Thomas Hartmann, Benjamin Zapilko, Joachim Wackerow, and Kai Eckert. Constraints to Validate RDF Data Quality on Common Vocabularies in the Social, Behavioral, and Economic Sciences. *Computing Research Repository (CoRR)*, abs/1504.04479, 2015. <http://arxiv.org/abs/1504.04479>.
8. Mary Vardigan, Pascal Heus, and Wendy Thomas. Data Documentation Initiative: Toward a Standard for the Social Sciences. *International Journal of Digital Curation*, 3(1):107 – 113, 2008. <http://www.ijdc.net/index.php/ijdc/article/view/66>.